

Deep Learning

Motive: - Basics, Maths, Interview Preparation

Day-1 Agenda

- 1) Deep learning \rightarrow Perception {AI vs ML vs DL vs DS}
- 2) Forward propagation
- 3) Backward Propagation
- 4) Loss function
- 5) Activation functions
- 6) Optimizers

Prerequisite

- ① Python
- ② ML
- ③ Stats

Q) Why deep learning is becoming popular?

\rightarrow Good Model Creation.

\rightarrow Hardware Advancement.

DL minimizes the need of human action since its algorithms conduct feature extraction on their own. This makes the process much faster and reduces the risk of human error.

Suppose:

$y=1$ (actual value)

$\hat{y}=0$ (predicted value)

then $y - \hat{y} = 1$ (this diff should be very near to 0)

↳ loss function

AI vs ML vs DL vs DS

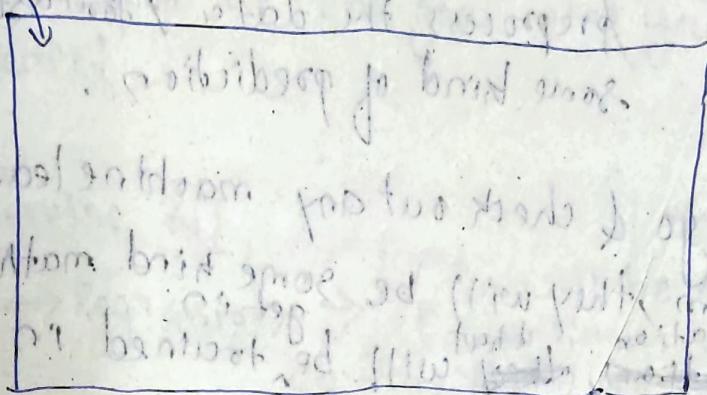
AI → Artificial Intelligence

ML → Machine Learning

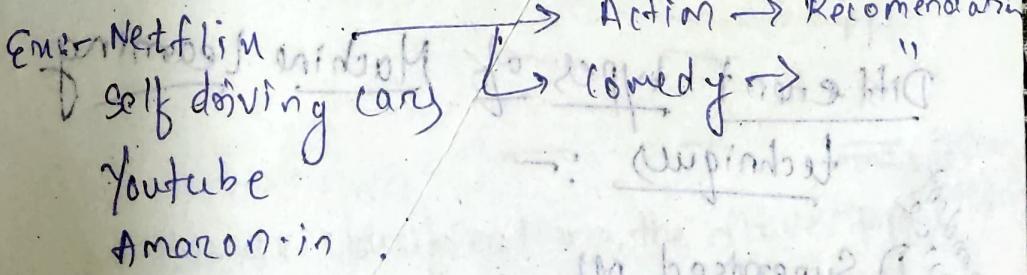
DL → Deep Learning

DS → Data Science

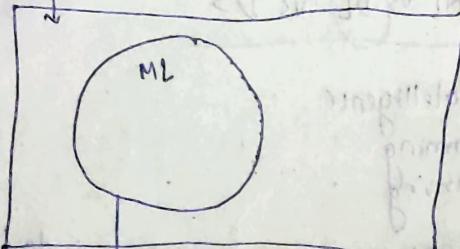
(It's an application which can do its own task without any human intervention)



At the end of the day we are creating an AI application.



AI → It means you have created an application and that application is performed without any human interference/intervention.



↳ Stats tools to analyze the data / visualize the data / preprocess the data / forecasting some kind of prediction.

If you go & check out any machine learning algorithm, they will be some kind mathematical equation that will be focused on the data.

Those mathematical equations can be considered as stats tools.

Different types of Machine Learning

techniques :-

Independent feature	Dependent feature	Weight	Output
Height	Height	130	110
Height	Height	140	120
Height	Height	150	130

It will train itself that whenever we give a height (input) it will give the specific weight (output)

In this ML Model : if we give height, it will be the independent feature and weight as the dependent feature (output feature). In short

→ In a supervise ML problem statement you will definitely have an output feature. So, you know us end goal what you need to predict.

1) Supervise ML

- Regression → O/P → will be a continuous value
Ex:- 45, 46, 47, 46.5, 46.3, -

→ Classification → O/P → fixed no. of categories

→ Independent feature → Height
→ Dependent feature → Gender
→ M
→ F

→ Binary classification

With binary classification, the above type will be binary classification as it will only have two possible outputs, i.e., 0 and 1.

If there will be more than two categories then it will be Multiclass classification.

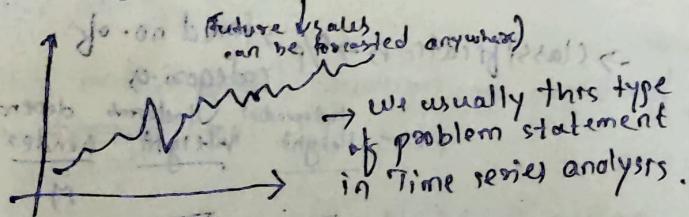
2) Multilabel Classification! -

(This multilabel classification is used in recommendation system as such)

Movies	<u>Action</u>	<u>Romance</u>	<u>Thriller</u>	<u>Suspense</u>
Endgame	1	0	1	1
Dhruvam	0	0	1	1

With respect to movies it can based on various categories, so based on probability we usually give us a/c/p.

Forecasting? → Regression problem



Multiclass classification

→ In multiclass, it has just one output only, and many no. of categories

Multilabel classification

→ In multilabel, it has many categories as well as many outputs.

2) Unsupervised ML :-

(clustering, segmentation, reduce dimension)

~~Ex:-~~

Company → product → costly product

Company has some kind of customer data such as
salary. $\begin{matrix} (0-9) \\ \text{Spending score(S)} \end{matrix}$ Email → product

65000

70000

75000

80000

100000

3

5

4

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

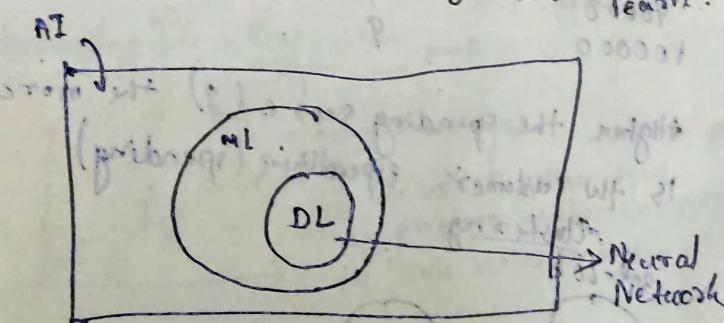
303

3) Semi Supervised ML :- (Recommendation)

In this ML,
we gonna have :

Labeled data → Model → Recommendation
with the interaction with mobile the film u watched such as action, romance

Deep Learning → Mimic human brain
→ make the machine learn.



1) ANN (Artificial Neural Network)

It works on all the tabular kind of data.

2) CNN (Convolutional Neural Network)

It works on images and videos.

- Applications:-
- 1) Image classification
 - 2) Object detection → RNN, LSTM, Bi-directional, PAPER CNN, GAN
 - 3) Object Segmentation
 - 4) Tracking
 - 5) GAN

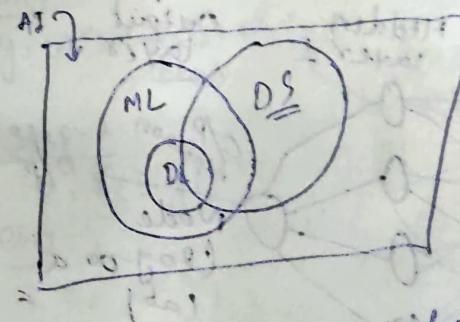
3) RNN (Recurrent Neural Network)

It works on text, time series and sequential data.

Techniques in RNN

- RNN
- LSTM RNN
- Bi-directional LSTM RNN
- Encoder-Decoder
- Transformers
- BERT
- GPT1, GPT2, GPT3

Data Science



NLP can be part of ML & DL.

NLP → Natural language Processing science

It is a subfield of linguistics, computer science

AI concerned with the interactions between computers and human language.

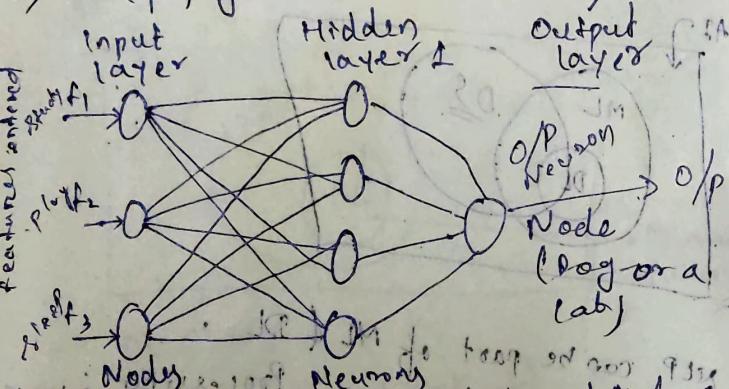
Deep Learning

- 1) ANN → (Artificial Neural Network)
- 2) CNN → (Convolutional Neural Network)
- 3) RNN → (Recurrent Neural Network)

It is a technique to mimic the brain's learning process.

Neural Network

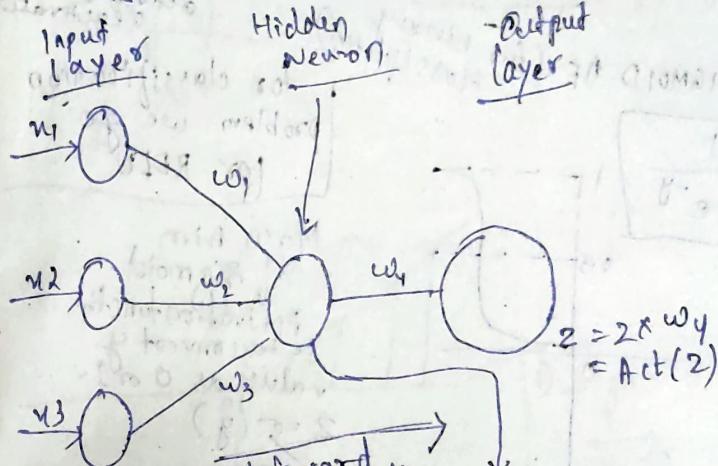
- Perceptron {Single layered Neural N/w}
- Back propagation (Jeffrey Hinton)



This is a simple Neural Network Architecture.

Dataset	ANN based			
	study	play	sleep	Pass/Fail
	7	3	7	1
	2	5	8	0
	4	3	7	1

How Neural Network Works



Step 1

$$y = w_1x_1 + w_2x_2 + w_3x_3 + \text{bias}$$

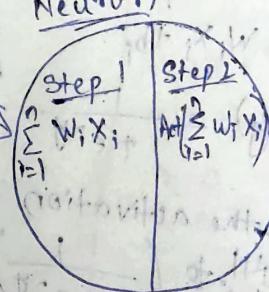
Step 2

$$z = \text{Act}(y)$$

Step 3

$$z = z + w_4$$

$$= \text{Act}(2)$$



$$w_1x_1 + w_2x_2 + w_3x_3$$

Bias → By mistake if you

add bias to x_1, x_2 and have initialized all

your weights as zero

bias will have some

value so that we

can continue the

behavior training

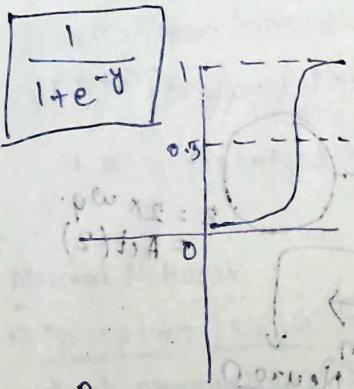
but this will be required

in every hidden layer,

there must be bias.

Activations Functions

1) SIGMOID AF (for binary classification)



$$y = \sum_{i=1}^n w_i x_i + b_i$$

$$\text{Act}(y) = \frac{1}{1 + e^{-y}}$$

Once the activation function is implemented it will be $\frac{1}{1 + e^{-y}}$ then the value

will be between 0 & 1.

If the value is less than 0.5, it will be considered as 0 and

If the value is more than 0.5, it will be considered as 1.

0.5 = threshold

O/p of sigmoid AF
= 0 or 1

0.5 ~~neuron isn't activated~~ less than 0.5

0.5 ~~neuron is activated~~ greater than 0.5

(whether the neurons should be activated or deactivated)

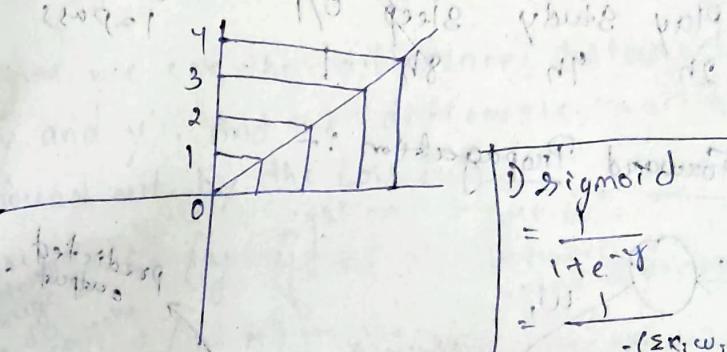
for classification problem we go for RELU

Main Aim of sigmoid Activation function is to convert y value as 0 or 1.
 $z = \sigma(y)$

for AF we solve non-linear problems.

2) RELU AF (for regression problem we go for RELU)

$$y = \sum_{i=1}^n w_i x_i + b;$$



1) Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

O/P \Rightarrow 0 or 1

$$\geq 0.5 = 1$$

$$< 0.5 = 0$$

If $y = -ve$

$$\max(-ve, 0) = 0$$

If $y = +ve$

$$\max(+ve, 0) = +ve$$

Importance of weights :-

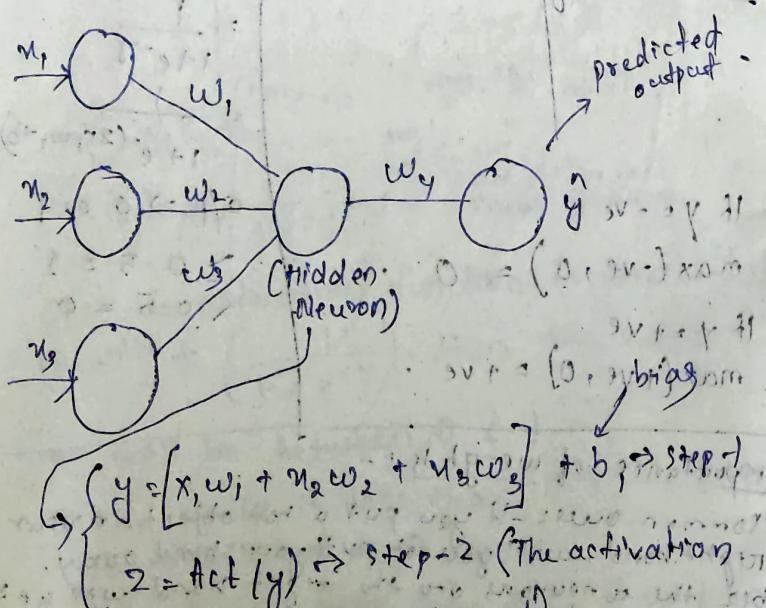
Common Ques - If you put a hot object in your right hand, will you not put your hand away, bcz the neurons in the right hand will get activated (i.e. Neuron will pass signal to brain saying that a hot object is placed and you have to remove your hand) so this weight will actually help your neurons to act.

Weights basically says that how much the neurons get activated or deactivated.

Neural Network Training with Back Propagation

We're having a dataset:			O = fail 1 = pass
(u ₁)	(u ₂)	(u ₃)	
Play	Study	Sleep	O/P
2h	4h	8h	1

Forward Propagation:



$$\{ y = [x, w_1 + u_2 w_2 + (u_3, w_3)] + b; \rightarrow \text{step-1} \}$$

$$\cdot 2 = \text{act}(y) \rightarrow \text{step-2} \quad (\text{The activation function may be sigmoid})$$

Sigmoid transfers the value betw. 0 - 1.

The average of all the loss values is called the error or loss function. This function tells us how far off our predicted output is from the actual output.

We'll discuss how to fit your predicted output to the actual output by minimizing the difference between them.

suppose $\hat{y} = 0.98$, then \hat{y} is real value
[and we already know that $y = 1$ (output)]
and we know, we need to compare this
that \hat{y} and y are almost same.
So, here we see the differences between
 \hat{y} and y . And the differences can be
found out by the Loss function.

If the Loss = $(\hat{y} - y)^2$ and if the difference
is huge, we do back propagation.
 $(1 - \hat{y})^2$ is also loss function.

Aint: To minimize the difference
so that \hat{y} & y is higher
if the diff is \uparrow we've to update
the weights which can be done through BP.

As we see the loss value is higher
i.e. 1 and by the difference of \hat{y} & y
it is completely & wrongly predicted.

So, in order to make it correct
we need to make changes in $w_1, w_2,$
 w_3 to make the $\hat{y} = 1$.

The loss value should be reduced
to a minimal value such that
 \hat{y} should match y .

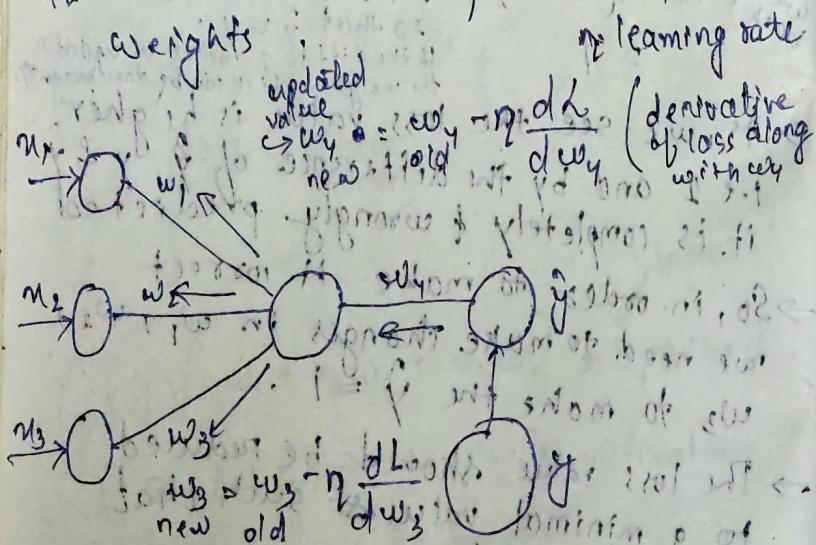
During the training phase we need
to change the w_1, w_2, w_3 to
decrease the loss values.

And it's can't be done with the help of
optimizers.

Eu of
Optimizers! - Gradient descent and some
(Suppose) gradient descent
with respect to Stochastic gradient descent
and many more (main aim of BP is)
Backward Propagation (to update weights)

As we got the loss value L .
Our optimizer will try to reduce the loss
value.
In order to reduce the loss value, we
basically have to back propagate now.

It means we have to update the pre
weights.



And learning rate should be a very
small value.

Let us consider $n = 0.001$.
The reason of taking n in a small value

because it will help us to go to the
global minimal function, of the
gradient descent whenever we want
to actually achieve the particular
point.

That is only possible by an optimizer.
Suppose we're considering gradient descent.

As hence, it has a single row of record.
Over here we derived a loss function.
If we have multiple rows of record
at that time we need to define the
cost function.

$$Cost = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Conclusion

FP

1) N/L layers

2) Weights

3) Bias

4) Activation function

Weight update formula:-

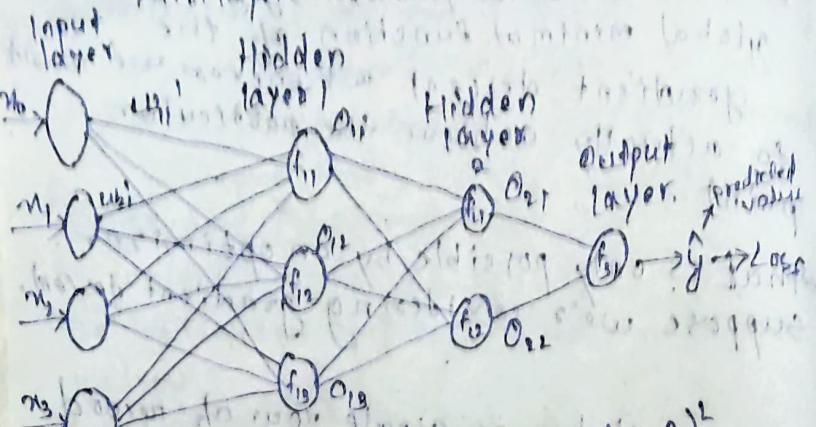
W_{new}

$$W_{\text{new}} = W_{\text{old}} - n \frac{dL}{dW_{\text{old}}}$$

$$\frac{dL}{dW_{\text{old}}}$$

- 5) Loss fu? ($\hat{y} - y$)
- 6) minimize the value
- 7) with the help of Optimizer
- 8) update the weights
- 9) update the bias also if needed.

Multilayer Neural Network



$$\text{Loss} = (y - \hat{y})^2$$

Input layer = 4x3

Optimizer \rightarrow Loss driven

HL1 \rightarrow 3x2
HL2 \rightarrow 2x1

algorithm used here is
gradient descent with learning rate

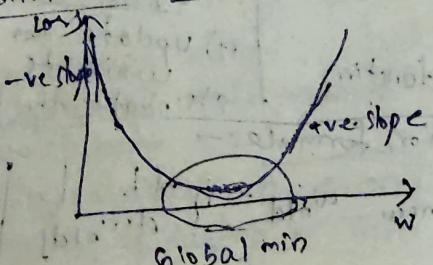
① Weight Updates

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{dL}{dW_{\text{old}}} \quad (B \cdot B)^T = 70$$

Q) why we found out derivatives?

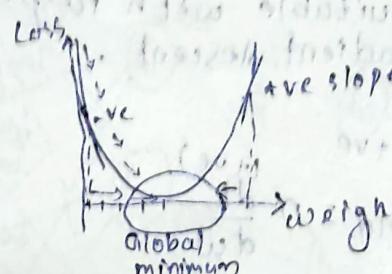
\rightarrow Derivative basically means we need to find out the slope.

GRADIENT DESCENT \rightarrow kind of steps:



Right hand side of the tangent is pointing down, then the slope is negative or vice versa

Derivatives will help us to find out whether this slope is positive or negative!



If slope = +ve (+ve)

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{dL}{dW_{\text{old}}}$$

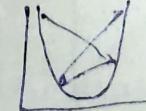
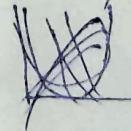
$$= W_{\text{old}} - (\eta \cdot +)$$

$$= W_{\text{old}} + (-)$$

Adding some values in older weight means points in graph is moving down towards the global minimum.

But the speed of going down depends on learning rate (η) .

If learning rate is small it will go slowly
If learning rate is very higher value then the points will jump from here and there and it may never reach the global minimal point.

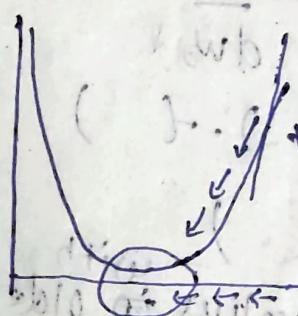


So, it is always says that we should select a learning rate in such a way that it should be a value which is suitable with respect to this gradient descent.

If slope :- +ve

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{dt}{dw}$$

$\hat{w}_{old} = \left(\hat{w}_{new} \right) \cdot \text{some value}$



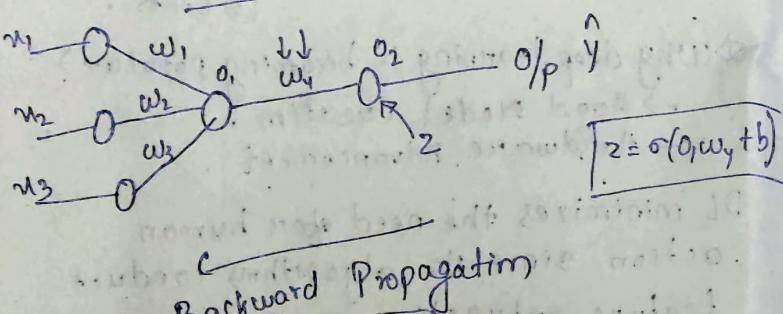
Day-2

Agenda

- 1) Forward Propagation → in front side
 - 2) Chain Rule of Derivatives
 - 3) Vanishing Gradient Problem
 - 4) Loss function
 - 5) Activation function

Chain Rule of Differentiation

$$L = L_{\text{loss}}$$



$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$$

w_{new} will be updated.

$$\rightarrow \text{How the weight } w_j \text{ changes?}$$

$\Rightarrow w_{j,\text{new}} = w_{j,\text{old}} - \eta \frac{dL}{dw_{j,\text{old}}}$

To find this value we particularly use chain rule of differentiation

chain rule of derivative

$$\frac{dL}{dw_{old}} = \frac{dL}{d\theta_2} * \frac{d\theta_2}{dw_{old}}$$

So, like weight activation function update in B.P
also bias activation function update

$$\checkmark b_{\text{new}} = b_{\text{old}} - \eta \boxed{\frac{dL}{db_{\text{old}}}}$$

Note :-

Note :-
Bias are basically the values that are added to the sums calculated at each node (except input nodes). b_0, w_0

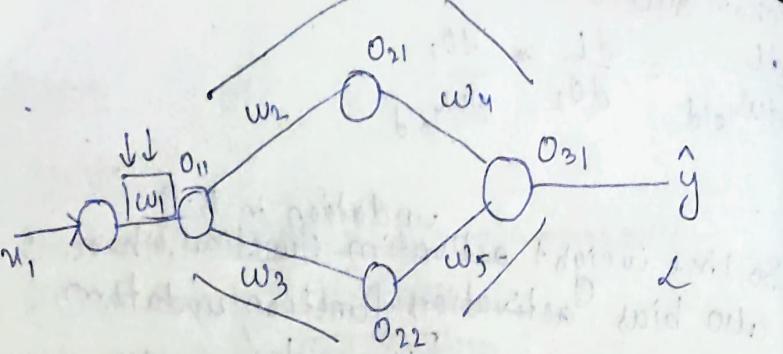
→ If we want to ~~find~~ update the weight of w_i ,

$$w_i^{\text{new}} = w_i^{\text{old}} - \eta \frac{\partial h}{\partial w_i} \quad \text{old}$$

$$\frac{dL}{d\omega_1}_{old} = \frac{dL}{d\omega_2} * \left(\frac{d\omega_2}{d\omega_1} \right)_{old} * \frac{d\omega_1}{d\omega_1}_{old}$$

$$\omega_{2\text{new}} = \omega_{2\text{old}} - \eta \frac{\omega_d L}{d\omega_{2\text{old}}}$$

$$\frac{dO_L}{d\omega_q} * \frac{d\omega_q}{dO_1}$$



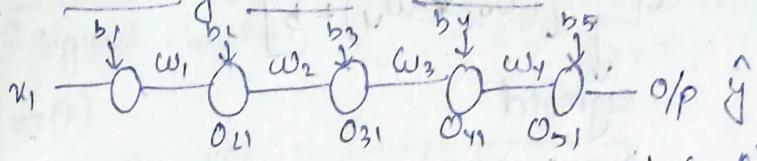
$$w_{1,\text{new}} = w_{1,\text{old}} - \eta \frac{dL}{dw_{1,\text{old}}}$$

$$\frac{dL}{dw_{1,\text{old}}} = \left[\frac{dL}{dO_{31}} * \frac{dO_{31}}{dO_{21}} * \frac{dO_{21}}{dO_{11}} * \frac{dO_{11}}{dw_{1,\text{old}}} \right] +$$

$$\left[\frac{dL}{dO_{31}} * \frac{dO_{31}}{dO_{22}} * \frac{dO_{22}}{dO_{11}} * \frac{dO_{11}}{dw_{1,\text{old}}} \right]$$

There are basically 2 ways to go from O_{31} to w_1 in B-P.
so, we add both the ways.

Vanishing Gradient Problem :-



$$(\text{Mean Squared Error}) \text{ MSE} = \text{Loss} = \frac{1}{2} (y - \hat{y})^2$$

$$w_{1,\text{new}} = w_{1,\text{old}} - \eta \frac{dL}{dw_{1,\text{old}}}$$

$$\frac{dL}{dw_{1,\text{old}}} = \frac{dL}{dO_{31}} * \frac{dO_{31}}{dO_{21}} * \frac{dO_{21}}{dO_{11}} * \frac{dO_{11}}{dw_{1,\text{old}}}$$

At every layer we use,

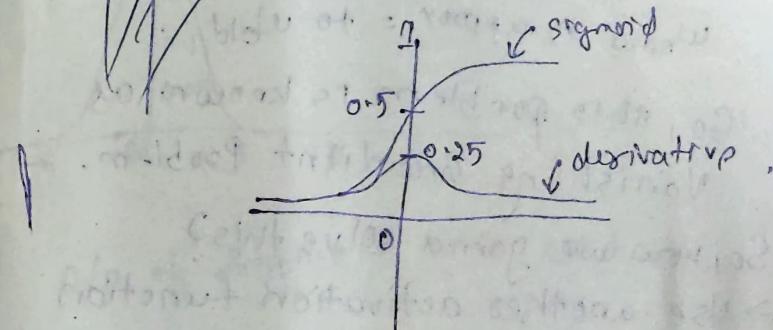
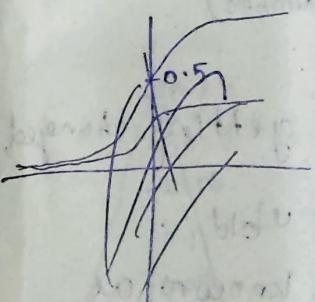
Sigmoid Activation fun

$$y = \frac{1}{1+e^{-x}}$$

O/p \rightarrow 0 or 1

≥ 0.5 1

< 0.5 0



$$\text{Derivative} = [0 \leq \sigma(y) \leq 0.25]$$

$$o_{51} = \sigma[(o_{41} * w_{41}) + b]$$

↓
sigmoid

$[0 \text{ to } 0.25]$

Let's take an example,
that

$$\frac{dL}{dO_{51}} + \frac{dO_{51}}{dO_{41}} * \frac{dO_{41}}{dO_{31}} + \frac{dO_{31}}{dO_{21}} * \frac{dO_{21}}{dw_1}$$

$$0.25 * 0.15 * 0.10 + 0.05 * 0.02$$

The values are getting ↓ (decreasing)
In short we're going to get a
very small value.

no
 $w_{\text{new}} = w_{\text{old}} - \eta \text{ (small number)}$

$$w_{\text{new}} \approx w_{\text{old}}$$

Hence w_{old} is hardly getting changed.

$$w_{\text{new}} \text{ is approx.} \approx w_{\text{old}}$$

So, this problem is known as
Vanishing Gradient Problem.

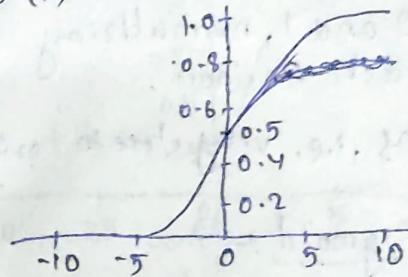
So, how we gonna solve this?

→ Use another activation function
instead of where there is no
vanishing gradient descent.

Activation functions

1) Sigmoid function

$$\sigma(x)$$



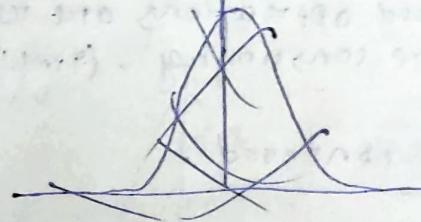
$$\sigma(u) = \frac{1}{1 + e^{-u}}$$

$$0.5 \geq u = 1$$

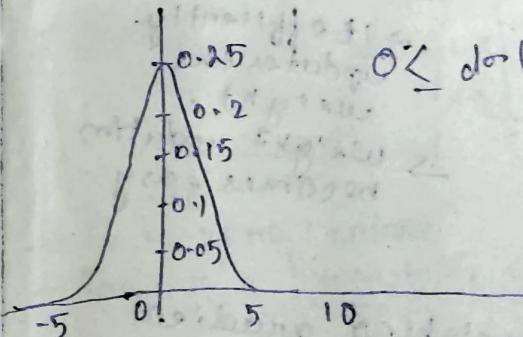
$$0.5 < u = 0$$

$$0 \leq \sigma(u) \leq 1$$

$$\frac{d\sigma(u)}{du} = \text{derivative}$$



$$0 \leq \frac{d\sigma(u)}{du} \leq 0.25$$



O/P of sigmoid = 0 or 1

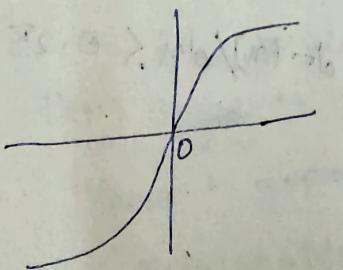
Advantages of Sigmoid :-

- 1) Smooth gradient, preventing "jumps" in output values
- 2) O/p bound bet' 0 and 1, normalizing the output of each neuron.
- 3) Clear Predictions, i.e. very close to 0 or 1.

Disadvantages of Sigmoid :-

- 1) Prone to gradient vanishing
- 2) function output is not zero centered
- 3) Power operations are relatively time consuming. (Time complexity ↑)

Zero-centered.



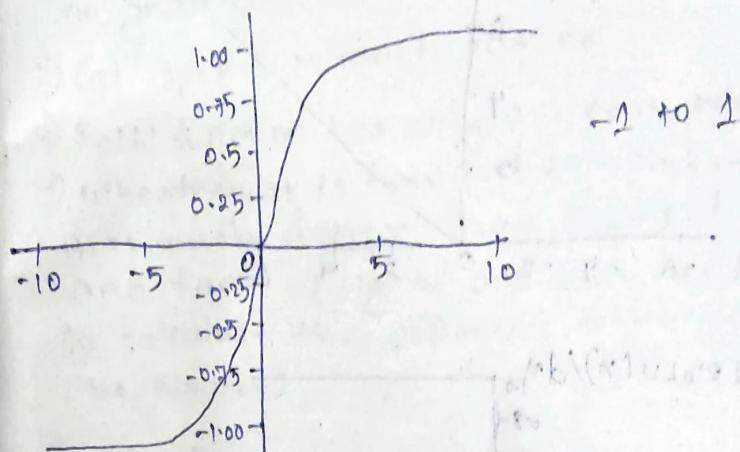
- IE efficiently updates the weight.
- Weight update becomes easy

If there is a vanishing gradient problem we can't create a deep neural network.
So, New activation fun' come up which are :-

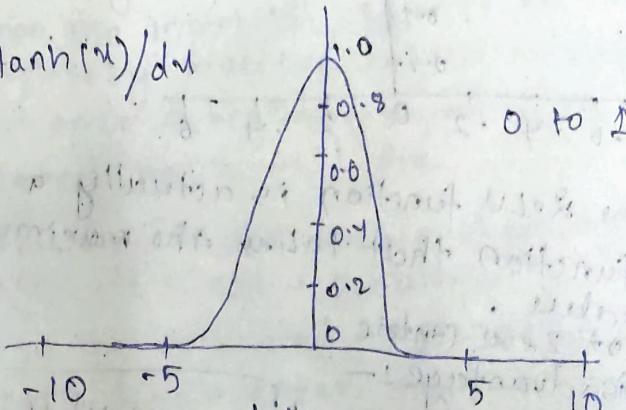
(Challenger's backprop)

Op tanh function.

$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



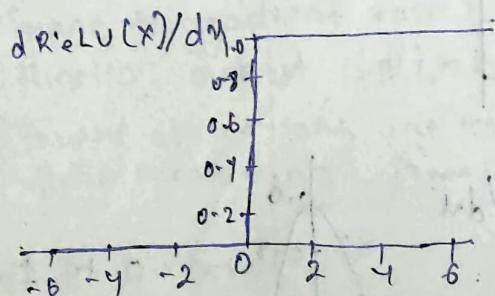
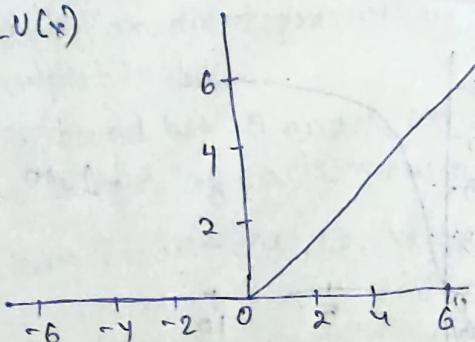
$d \tanh(x)/dx$



- It is zero centered
- Tanh is a hyperbolic tangent function
- It is better than sigmoid activation function in case of vanishing gradient descent.
- But if we create a very deep neural network, vanishing gradient problem will still exist.

3) ReLU function (Rectified Linear Unit)
 $\text{ReLU}(x) = \max(0, x)$

$\text{ReLU}(x)$



→ The ReLU function is actually a function that takes the maximum value.

→ Not zero centric

Disadvantage:-

$$\frac{dL}{dW_{old}} = 0$$

$W_{new} \neq W_{old}$

$\cancel{x} \rightarrow \text{Dead}$

Dead
state

0

0

Advantages:-

1) When the input is positive, there is no gradient saturation problem.

2) Calc' speed is much faster

3) ReLU function has only a linear relationship

4) Whether it is forward or backward, it is much faster than sigmoid and tanh. (Sigmoid and tanh need to calculate the exponent, which will be slower)

Disadvantages:-

1) When the input is negative, ReLU is completely inactive. which means that once a negative number is entered, ReLU will die.

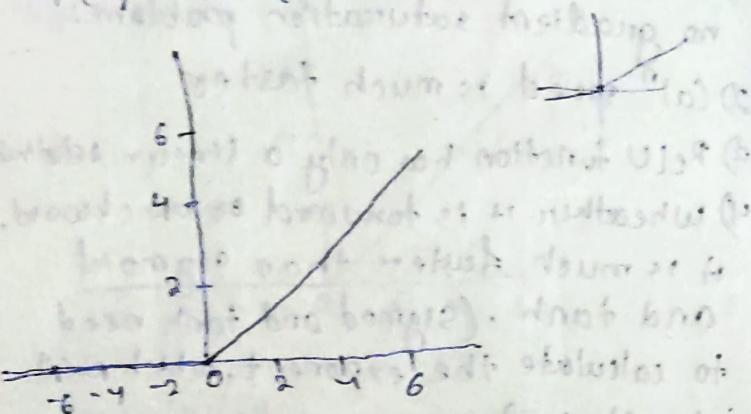
In this way, in the forward propagation process, it is not a problem.

But in the backward propagation process, if you enter ~~a negative~~ a negative number, the gradient will be completely zero, which has the same problem as the sigmoid function and tanh function.

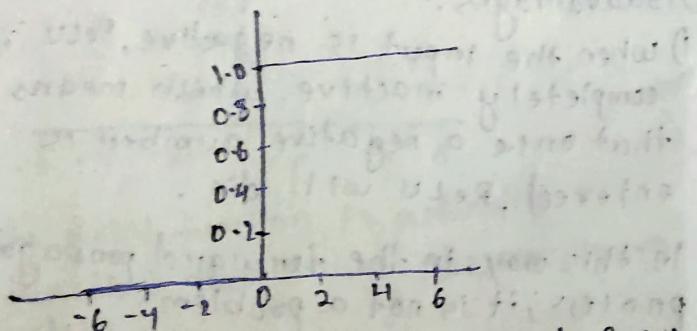
2) We find that the output of the ReLU function is either 0 or a positive number, which means that the ReLU function is not a 0-centrict function.

→ Leaky ReLU function (Solve Dead Neuron)

$$f(u) = \max(0.01u, u)$$



$$\frac{df(u)}{du} = \text{Derivative}$$

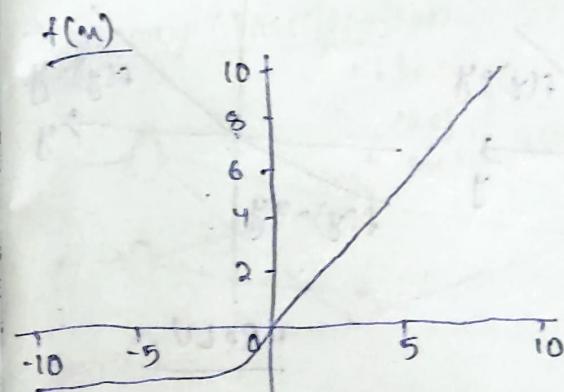


→ In order to solve the Dead ReLU problem, we use Leaky ReLU to set the first half of ReLU 0.01u instead of 0.

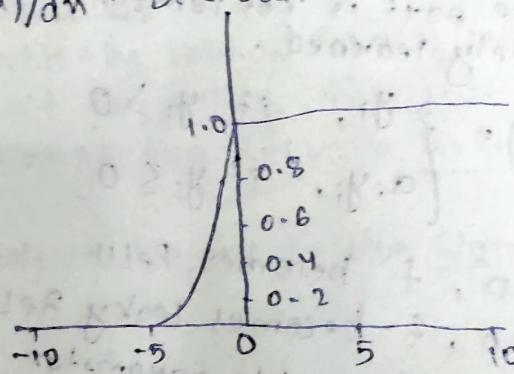
→ In theory, Leaky ReLU has all the advantages of ReLU, plus there will be no problems with Dead ReLU, but in actual operation, it hasn't been fully proved that Leaky ReLU is better than ReLU.

→ ELU (Exponential Linear Units) function

$$f(u) = \begin{cases} u, & \text{if } u > 0 \\ \alpha(e^u - 1), & \text{otherwise} \end{cases}$$



$$\frac{df(u)}{du} = \text{Derivative}$$



→ ELU is also proposed to solve the problems of ReLU.

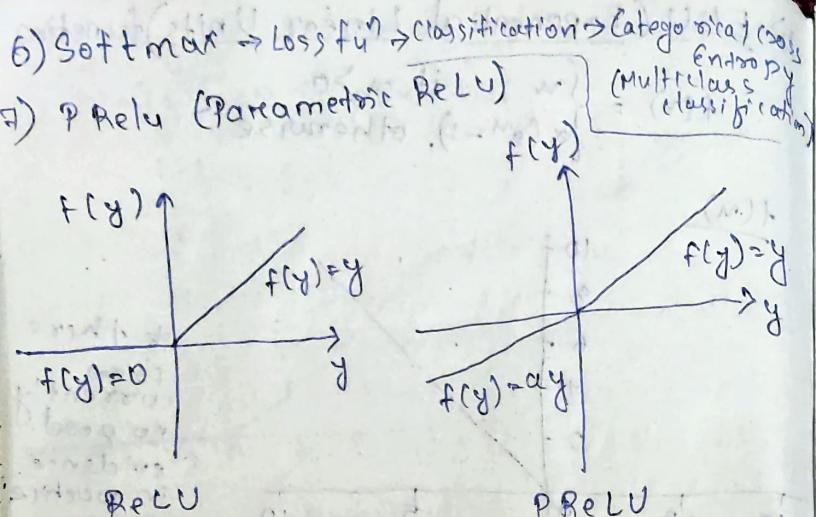
→ ELU has advantages over ReLU

→ No dead ReLU issues

→ the mean of the output is close to 0, zero-centered.

Disadvantage: - It is slightly more computationally intensive. Similar to leaky ReLU, although theoretically better than ReLU.

at there
is very
currently
no good
evidence
in practice
that ELU
is always
better
than ReLU



For PReLU, the coefficient of the negative part is not constant and is adaptively learned.

$$f(y_i) = \begin{cases} y_i, & \text{if } y_i > 0 \\ a_i y_i, & \text{if } y_i \leq 0 \end{cases}$$

If $a_i = 0$, f becomes ReLU

If $a_i > 0$, f becomes leaky ReLU

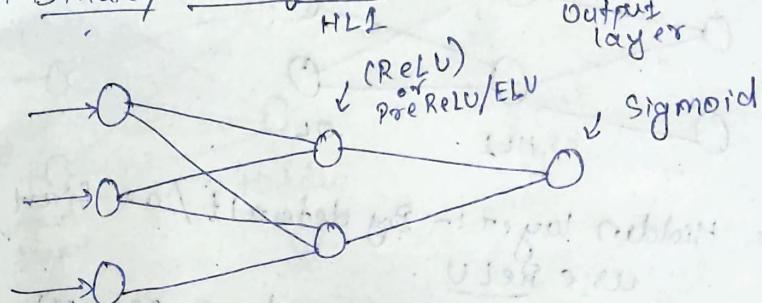
If a_i is a learnable parameter,

f becomes PReLU

We have to put a_i as learnable parameter.
it keeps on changing the parameters.

Technique which Activation fn we should use

Binary Classification

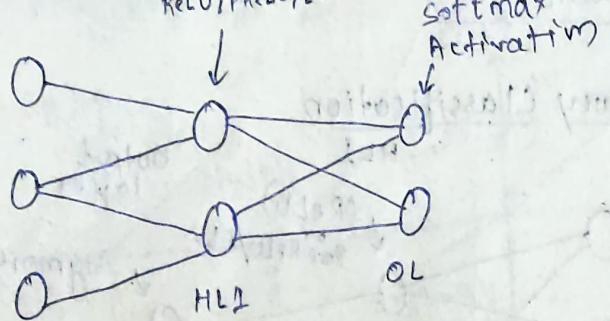


* Hidden layer :- Use ReLU at first.
~~if ReLU isn't~~
If with the help of ReLU if convergence isn't happening, then we can change the ReLU to PreReLU or ELU

* Output Layer :- Use Sigmoid

Multiclass Classification

ReLU/PReLU/ELU

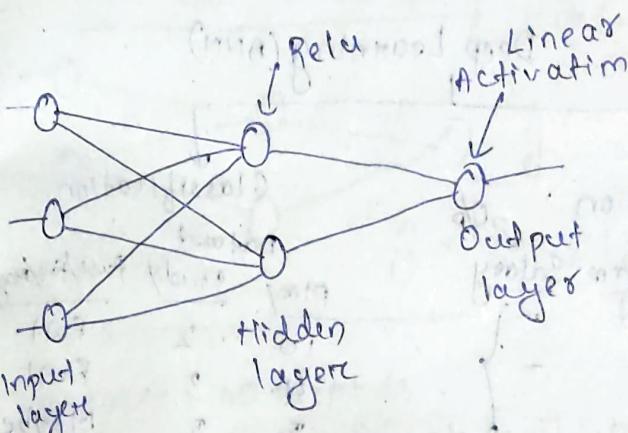


* Hidden layers:- By default / at first use ReLU.

If dead state situation occurs then use PReLU or ELU

* Output layer:- Use softmax

Regression (It means u gonna have continuous value)



* Hidden layers:- Use ReLU or any variation of ReLU.

* Output layer :- Linear Activation function and over here there will be a separate loss function

Loss Functions

Deep Learning (ANN)

Regression

Exp	Degree	Salary
10	PhD	100
-	-	-
-	-	-
-	-	-

Regression Problem

Classification

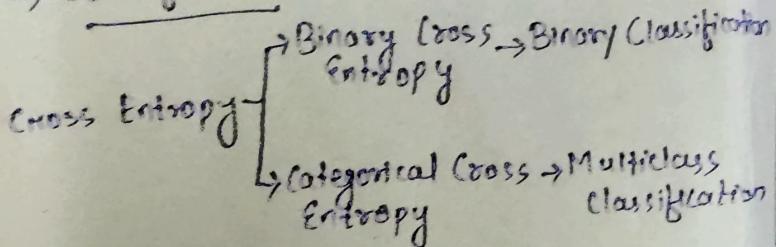
Dataset	Play	Study	Pass/Fail
10	2	Fail	
4	3	Fail	
5	3	Maybe	
2	7	Pass	

This dataset is a multiclass classification

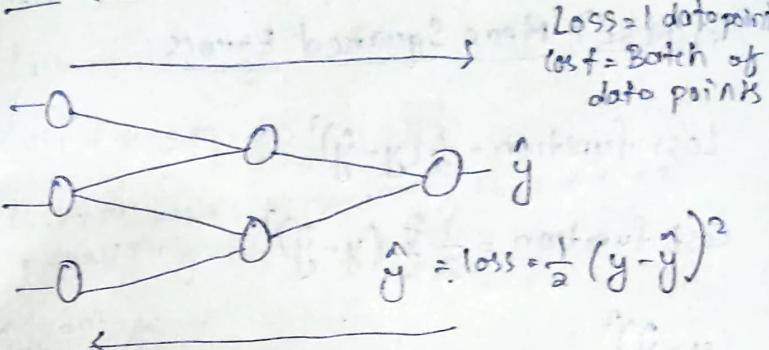
Regression

- MSE (Mean Squared Error)
- MAE (Mean Absolute Error)
- Huber Loss

Classification



Loss Function AND Cost Function



Loss = 1 datapoint
Cost f = Batch of data points

Dataset = 100 records.

→ This loss function gets calculated for every record.

→ But this is not an efficient scenario. Suppose if I have 100 records I can set up a batch size. I can basically say at every forward propagation pass 10 records.

When you are passing the 10 records then you calculate the loss for the 10 records.

Whenever I pass batch at once. Then the formula get changes instead of loss I will try to write this as cost function.

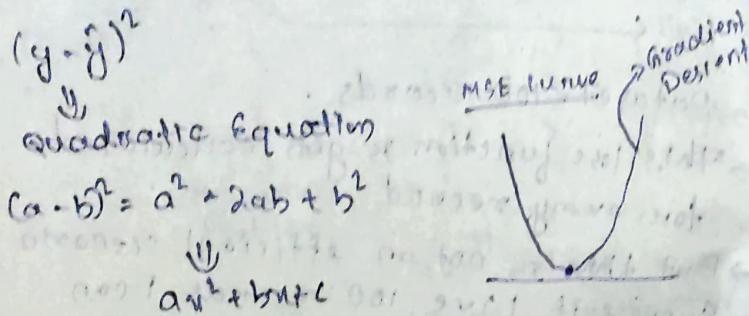
$$\text{Cost} = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Regression

i) (MSE) Mean Squared Error

$$\text{Loss function} = \frac{1}{2} (y - \hat{y})^2$$

$$\text{Cost function} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



Advantages

- 1) Differentiable
- 2) It has only 1 local or global minimum.
- 3) It converges faster

Disadvantage

- 1) Not Robust to outliers

The reason of shifting the line so much is because we're penalising the error. (squaring the error). Because of that the best fit line will have major shift

ii) (MAE) Mean Absolute Error

$$\text{Loss function} = \frac{1}{2} |y - \hat{y}|$$

$$\text{Cost function} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

As, here we're not squaring it, we're not penalising it.

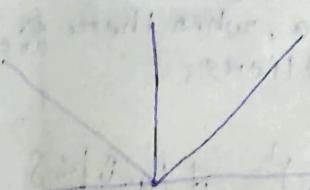
Advantages

- 1) Robust to outliers

Hence there will be a minor shift not a major one as of MSE.



MAE curve



→ It is time consuming

Hence the derivative will not be simple to calculate. We've to take subgradient.

Subgradient: We've to divide this fine part by part and calculate the slope. (Bcz here there is no curve)

iii) Huber Loss

→ It is a combination of MSE and NAE.

$$\text{Loss} = \begin{cases} \frac{1}{2}(y - \hat{y})^2, & \text{if } |y - \hat{y}| \leq \delta \\ \delta|y - \hat{y}| - \frac{1}{2}\delta^2, & \text{otherwise} \end{cases}$$

δ = hyperparameter.

$\boxed{\frac{1}{2}(y - \hat{y})^2}$ → we will use this formula

when there is no outliers:

⇒ If the error difference $|y - \hat{y}|$ is less than the hyperparameter.

$\boxed{\delta|y - \hat{y}| - \frac{1}{2}\delta^2}$ → we will use this formula, when there are ~~are~~ outliers.

$$\text{Cost function} = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \text{ if } |y_i - \hat{y}_i| \leq \delta$$

2) Classification

i) Binary Cross Entropy (Binary Classification)

$$\text{Loss} = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y})$$

Log loss function.

This is the loss that is used in logistic Regression.

As logistic regression is a binary classification problem.

We can write the above formula as :-

$$\text{Loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

$$\hat{y} = \frac{1}{1+e^{-x}} \quad (\text{sigmoid activation function in the last layer})$$

i) Categorical Cross Entropy (Multiclass Classification)

f_1	f_2	f_3	O/p
2	3	4	Good
5	6	7	Bad
8	9	10	Neutral

Good Bad Neutral

1	0	0
0	1	0
0	0	1

$$\text{Loss} = L(u_i, y_i) = - \sum_{j=1}^c y_{ij} * \ln(\hat{y}_{ij})$$

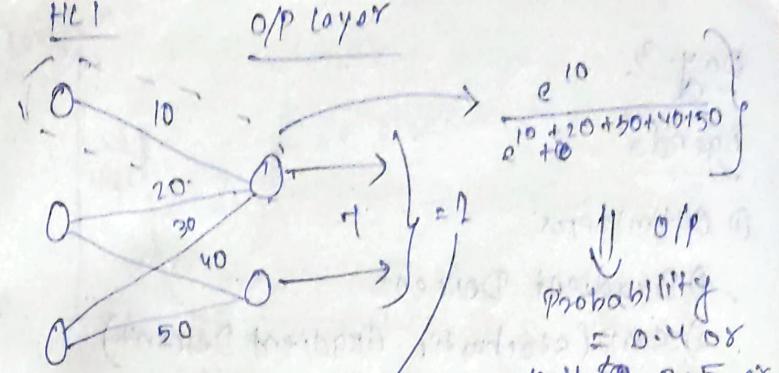
$$y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}] \quad c = \text{no. of categories}$$

$$y_{ij} = \begin{cases} 1 & \text{if the element is in class } j \\ 0 & \text{otherwise.} \end{cases} \quad i = \text{rows}$$

$$\hat{y}_{ij} = \text{Softmax Activation} \quad (\text{APPLIED in the O/p of the multiple classifiers})$$

Softmax Activation formula :-

$$\sigma(z) = \frac{e^{z_j}}{\sum_{j=1}^c e^{z_j}}$$



$$\text{Probability} = 0.4 \text{ or } 0.5 \text{ or } 0.6$$

Adding both the probabilities from the output layer will be ≈ 1 : whichever have higher probability, that class output will be coming up.

* Softmax Activation Function is applied on every output layer for a multiclass classification problem.

i) Categorical Cross Entropy (Multiclass Classification)

	f_1	f_2	f_3	O/P
1	3	4	5	Good
2	6	7	8	Bad
3	9	10	11	Neutral

Good Bad Neutral

1	0	0
0	1	0
0	0	1

Loss =

$$L(u_i, y_i) = - \sum_{j=1}^c y_{ij} * \ln(\hat{y}_{ij})$$

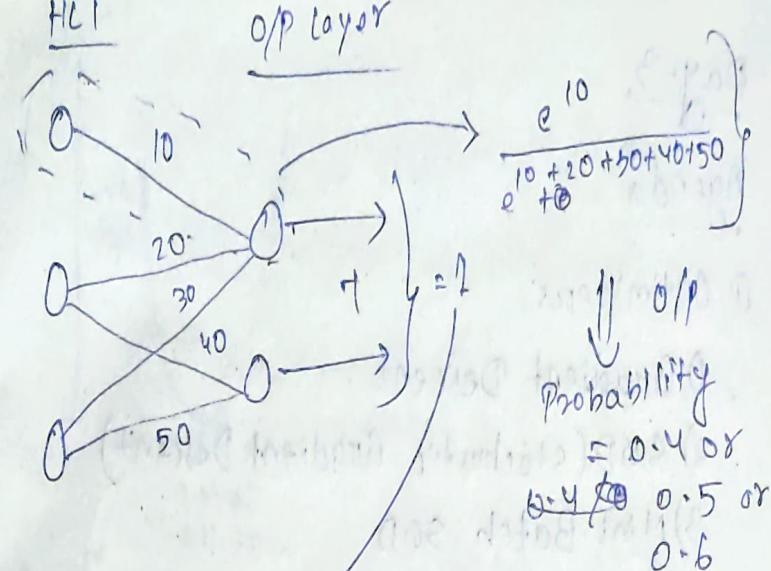
$$y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}] \quad \begin{matrix} c = \text{no. of categories} \\ i = \text{rows} \\ j = \text{column} \end{matrix}$$

$$y_{ij} = \begin{cases} 1 & \text{if the element is in class} \\ 0 & \text{otherwise} \end{cases}$$

\hat{y}_{ij} = Softmax Activation (Applied in the O/p of the multiple classification)

Softmax Activation formula :-

$$\sigma(z) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$



Adding both the probabilities from the output layer will be = 1: whichever have higher probability, that class output will be coming up.

* Softmax Activation Function is applied on every output layer for a multiclass classification problem.

Day-3.

Agenda

① Optimizers

- 1) Gradient Descent
- 2) SGD (Stochastic Gradient Descent)
- 3) Mini Batch SGD
- 4) SGD with Momentum
- 5) Adagrad
- 6) RMSprop
- 7) Adam Optimizers

- ⑧ Batch \rightarrow How many records to be given in a single Epoch
 (Total no. of times the algo will go through the entire dataset)
- ⑨ Epochs \rightarrow 1 F.P + 1 B.P = 1 Epoch
- ⑩ Iterations \rightarrow total no. of batches needed to complete 1 epoch

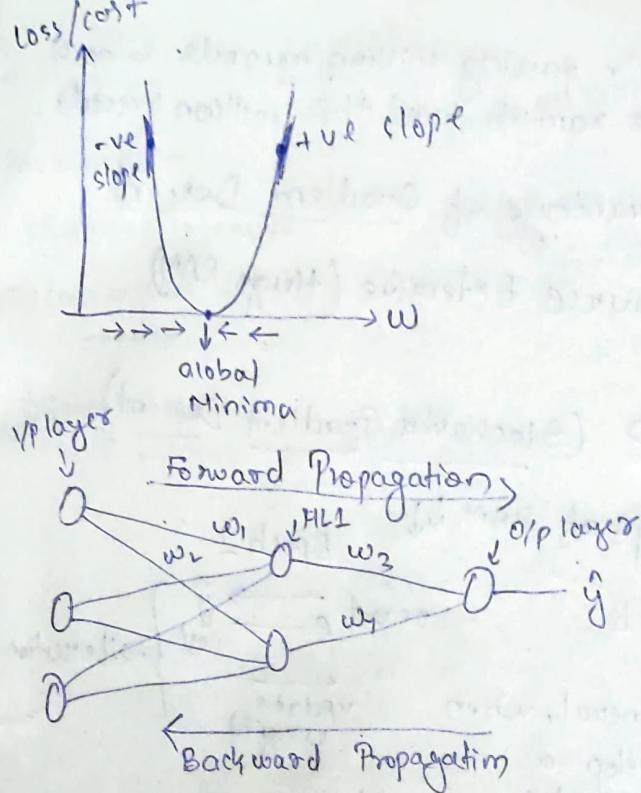
① Optimizers :-

Gradient Descent

Weight Updation Formula

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$$

learning rate.



$$\text{Cost function} = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$\left. \begin{array}{l} \text{1 Forwarded Propagation} \\ + \\ \text{1 Backward Propagation} \end{array} \right\} = 1 \text{ Epoch}$

If you have million records in the dataset in the F.P at a batch you're taking million records and finding out the cost function so we'll $\sum_{i=1}^n$ million.

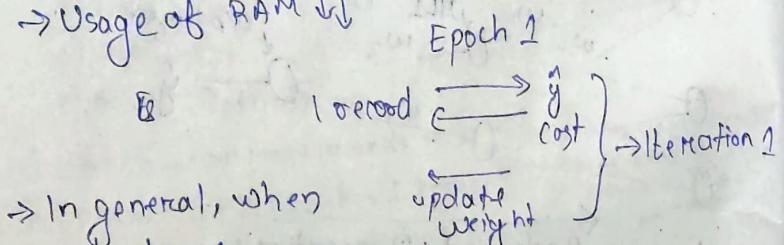
So, if you're passing million records you need a huge ram to load this million records

Disadvantage of Gradient Descent

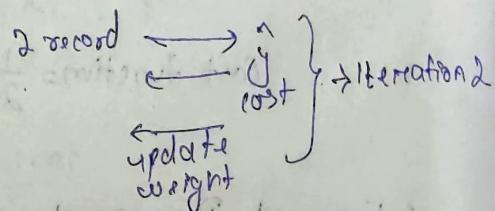
1) Resource Extensive (Huge RAM)

2) SGD (Stochastic Gradient Descent)

→ Usage of RAM ↓



→ In general, when we develop a deep learning model, there will be more than 1 epoch.



Suppose let's take 100 epoch

So, 100 million iterations.
It will make the progress slow.

} million iterations

Disadvantage of SGD

- 1) Convergence will be very slow.
- 2) Time complexity will be very high

3 Mini Batch SGD

1000000 → records

Advantage:-

1) Resource Intensive

2) Convergence will be better

batch size = 1000

$$\text{Total no. of iterations} = \frac{1000000}{1000} = 1000$$

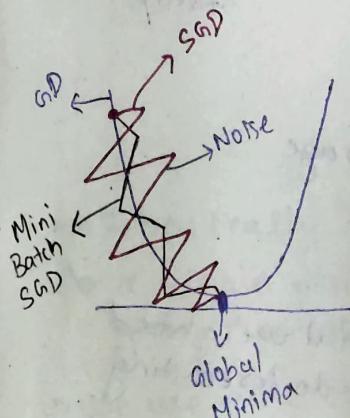
3) Time Complexity will improve

epoch-1

1000 } Iteration 1

→ } Iteration 2

→ } Iteration 1000



Noise:- The deviation from the path.

SGD > MiniBatch > GD
SGD

Q) How do we remove the noise?
→ We use momentum

* Mini batch SGD gives less noise as compared to SGD. But there is noise to remove it we use mini batch SGD with momentum.

4
y
Mini Batch SGD

SGD with Momentum

Momentum will smoothen the journey and reduce the noise & try to come to the global minima.

Q: How will it smoothen?
→ With the help of exponential moving average

Exponential Weighted / Moving Average



Time Series



ARIMA, ARMA (Models where EWMA is used)

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$$

t = current time stamp

$$w_t = w_{t-1} - \eta \frac{dL}{dw_{t-1}}$$

Exponential Weighted Average

$$t_1 \ t_2 \ t_3 \ t_4 \dots \ t_n$$

$$a_1 \ a_2 \ a_3 \ a_4 \dots \ a_n$$

With the help of exponential weighted average we can also do forecasting.

$$v_{t+1} = a_1$$

$$v_{t+2} = \beta * v_{t+1} + (1 - \beta) * a_2$$

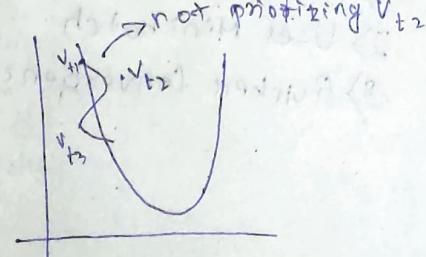
β = hyperparameter

β = It helps to decide which point to prioritize more, prev. one or the present one.

Suppose $\beta = 0.95$

$$\begin{aligned} v_{t+2} &= (0.95) * v_{t+1} + (1 - 0.95) * a_2 \\ &= 0.95 * v_{t+1} + (0.05) * a_2 \end{aligned}$$

Here, we can see that it's prioritizing the previous point (v_{t+1})



→ It basically gives the control to a single point betⁿ the two given points.

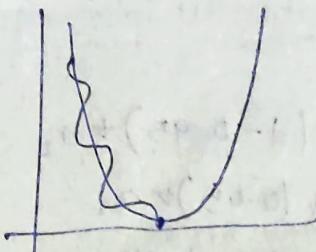
→ We use the exponential weighted average in the derivative.

$$w_t = w_{t-1} - \eta \frac{dL}{dw_t}$$

$$\Rightarrow w_t = w_{t-1} - \eta V_{dw}$$

$$V_{dw_t} = \beta \times V_{dw_{t-1}} + (1-\beta) * \frac{dh}{dw_{t-1}}$$

\downarrow
Exponential weighted Avg.



Advantages:-

- 1) Reduces the Noise
- 2) Uses Minibatch
- 3) Quicker Convergence

2 Adagrad (Adaptive Gradient Descent)

Objective:- As we've a learning rate in weight updation. And it is fixed till now. Can we do a certain thing that :-

At first, when the point is far away from global minima we can increase the speed by increasing the learning rate.

Later, when the point is near the global minima we can decrease the learning rate.

Can we do this?

→ Yes in Adagrad.

$$w_t = w_{t-1} - \eta \frac{dL}{dw_{t-1}} \quad \eta = \text{fixed.}$$

$$\boxed{w_t = w_{t-1} - \eta' \frac{dL}{dw_{t-1}}} \quad \eta' = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$$

$$\alpha_t = \sum_{i=1}^t \left(\frac{dh}{dw_i} \right)^2$$

As $\alpha_t \uparrow$ $\eta' \downarrow$ $\because \alpha_t$ divides η which is fixed

Bcz of this, we can see that as we go near the global minima, η' decreases.
 $\eta = 0.01$

ϵ is added to avoid divide by zero

By this way we're bringing adaptiveness to the learning rate.

η = fixed \Rightarrow adaptive \Rightarrow learning Rate =

Decreasing \Rightarrow Global Minima

$$t=1 \quad \eta = 0.01 \quad t=2 \quad \eta = 0.005 \quad t=3 \quad \eta = 0.002$$

Here η' will be a very small no. then $w_t \approx w_{t-1}$

& Adamelta And RMSprop

$$\eta' = \frac{\eta}{\sqrt{sdw + \epsilon}}$$

Exponential Weighted Average:-

$$sdw = \beta sdw_{t-1} + (1-\beta) \left(\frac{dw}{dw_{t-1}} \right)^2$$

$$\text{Suppose } \beta = 0.95$$

$$sdw_t = (0.95) sdw_{t-1} + (0.05) \left(\frac{dL}{dw_{t-1}} \right)^2$$

This will not increase that much

$$\Rightarrow w_t = w_{t-1} - \eta' \frac{dL}{dw_{t-1}}$$

$$\Rightarrow w_t = w_{t-1} - \frac{\eta}{\sqrt{sdw + \epsilon}} \frac{dL}{dw_{t-1}}$$

Combining both $\frac{sdw}{\sqrt{sdw + \epsilon}} + \frac{Vdw}{\sqrt{sdw + \epsilon}} = \text{Adam}$

low rating the momentum
 η at a point

Adam Optimizer (Best Optimizer)

Momentum + RMSprop (Adaptive learning Rate)
 $[Vdw] \quad [sdw]$

Adam Optimizer

$$Vdw = 0 \quad Vdb = 0 \quad sdw = 0 \quad sdb = 0$$

$$w_t = w_{t-1} - \eta' Vdw$$

$$b_t = b_{t-1} - \eta' Vdb$$

$$\eta' = \frac{\eta}{\sqrt{sdw + \epsilon}}$$

$$Vdw_t = \beta \times Vdw_{t-1} + (1-\beta) \frac{dL}{dw_{t-1}}$$

$$Vdb_t = \beta \times Vdb_{t-1} + (1-\beta) \frac{dL}{db_{t-1}}$$

ϵ = small no
 ϵ is added to avoid divide by zero

Day-4

Agenda

- 1) ANN practical implementation
- 2) Early stopping
- 3) Black Box Models vs white Box Models
- 4) CNN introduction

Black Box Models Vs white Box Models

Random Forest → Black Box Model

Decision Tree → White Box Model

Linear Reg → White Box Model

ANN → Black Box Model

Xgboost → Black Box Model

Black Box → We can't see how the algo works as it's quite huge and quite big internally.

white Box → we can't observe the internal

1. Non Tracing
2. Not Possible

Day - 5

CNN

Convolution Neural Network

(Image, Video frame)

Agenda

1) CNN vs Human Brain

2) Convolutional Operation

 → Convolution

 → Padding

 → Stride

 → Filters (kernels)

3) Max Pooling

4) Flattening

5) Practical Implementation.

CNN vs Human Brain

CAT	DOG
BICYCLE	
HUMAN	CAR

CEREBRAL CORTEX



VISUAL CORTEX

V₁

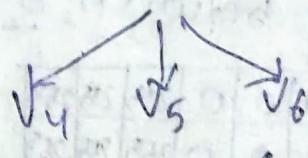
→ Having objects

V₂

→ Animals (CAT)

V₃

→ MAP the environment



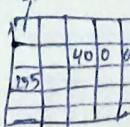
V₇

OP

2 Convolution Operation

Convolution

$0 \sim 255$

Images =  5×5 pixel

'Black & white' &
'RGB'

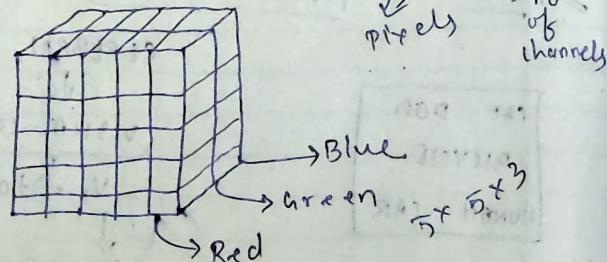
→ We divide the image into some kinds of pixels.

→ Each pixel range = $0 \sim 255$

Left $\rightarrow 0 = \text{Black}$
 $\rightarrow 255 = \text{White}$

→ Every pixel may have different values based on the image composition.

RGB (Coloured Image)



→ All 3 combines to form a image.

→ Here also the values will be ranging between 0 to 255 .

0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255

\otimes^6

→ We divide each pixel with 255

→ To get O/P as 0 or 1

⇒ Step-1 in convolution = bring all the values before

$0 \sim 1$

Placing

0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1

convolution operation

3×3

Kernel/filter
values:
horizontal edge filter

1	2	1
0	0	0
-1	-2	-1

0	0
0	0
0	0
0	0
0	0

4×4
O/P

→ Dividing each pixel with 255 is called as mix-max scaling.

→ The values of the 3×3 kernel/filter is known as horizontal edge filter.

→ Step-2:- Placing the filter in the image.

→ Then we perform multiplication of each shell of the image with each shell of the filter.

→ Then add all of it and place it in the O/p (4×4).

→ Calculation for the 1st instance:-

$$= (0 \times 1) + (0 \times 2) + (0 \times 1) + (0 \times 0) + (0 \times 0) + (0 \times 0) + (0 \times -1) + (0 \times -2) + (0 \times -1)$$

= 0 (It will be stored in the 1st shell of the 4×4 O/P shift)

→ Then move the filter one place towards the right.

Stride = Shifting/Moving the filter towards right by a certain value.

Stride = 2

↳ shifting by 1

Stride = 2

↳ shifting the filter in the image by 2.

Here, we've taken stride = 1

→ Calculation for the 2nd instance :-

$$\begin{aligned}
 &= (0x1) + (0x2) + (1x1) + (0x0) + (0x0) \\
 &\quad + (0x-1) + (0x-2) + (1x-1) \\
 &= 0 + 0 + 1 + 0 + 0 + 0 + 0 - 1 \\
 &= 0 \text{ (It will be stored in the 2nd cell)}
 \end{aligned}$$

→ 3rd & 4th instance = 0
 → After it completes the shifting towards right, shift to down, from the starting point of that row.
 Shift to the down by 1

→ So, the final O/P ↗

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

0	0	0	0
-4	4	-4	-4
-4	4	-4	-4
0	0	0	0

→ The above operation is specifically called as convolution.

0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1

(vertical edge filter)

1	0	-1
2	0	-2
1	0	-1

0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-1	-1	0

Min-Max Scaling :-

Smaller Value → 0
 Higher Value = 255

will become

Applying min-max scaling

255	0	0	255
255	0	0	255
255	0	0	255
255	0	0	255

smaller value

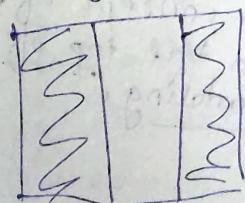
$$= -4 \rightarrow 0$$

higher value

$$= 0 \rightarrow 255$$

0 → black colour

255 → white colour



Vertical edge

Filter → Round shape filter
Horizontal
Vertical

$$\begin{array}{c} \text{(n-f) } \\ \text{size of image} \\ \underline{6 \times 6} \end{array} \quad \begin{array}{c} \text{(f=3) filter} \\ \underline{3 \times 3} \end{array} \quad \begin{array}{c} \text{O/P} \\ \Rightarrow 4 \times 4 \end{array}$$
$$[n-f+1] = 6-3+1 = 4$$

Here, the image size is decreasing which is losing some information.

Here the O/P image size is decreasing from the original I/P image.

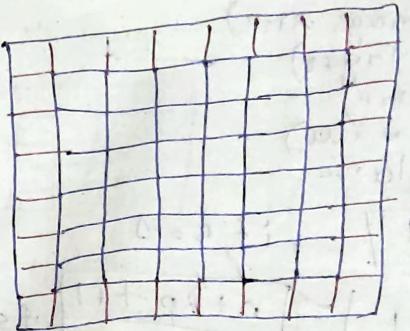
→ which means losing some information.

→ To avoid this, we can come up with something called padding.

→ If you've a image of 6×6 & wants an O/P of image of 6×6 . Then apply padding on the original I/P image.

→ The concept of applying a layer on the top of the image is called padding.

→ ~~Info~~



After applying padding on a 6×6 image.

It will become a 8×8 image.
(But the actual image will be only 6×6)

There will be an external layer added to it.

Types of padding:-

1) zero padding

2) whatever the nearest value to the padding, put that value into padding.

Before Padding :-

$$\begin{array}{l} n=6 \\ f=3 \end{array} \quad [n-f+1]$$

$$\rightarrow 6-3+1 \Rightarrow 4 = \text{O/P}$$

Here, so image size is decreasing, information loss takes place.

After padding :-

$$\begin{array}{l} n=8 \\ f=3 \end{array} \quad [n-f+1]$$

$$\rightarrow 8-3+1 \Rightarrow 6 = \text{O/P}$$

Here, so actual image size is maintained. No. information is lost.

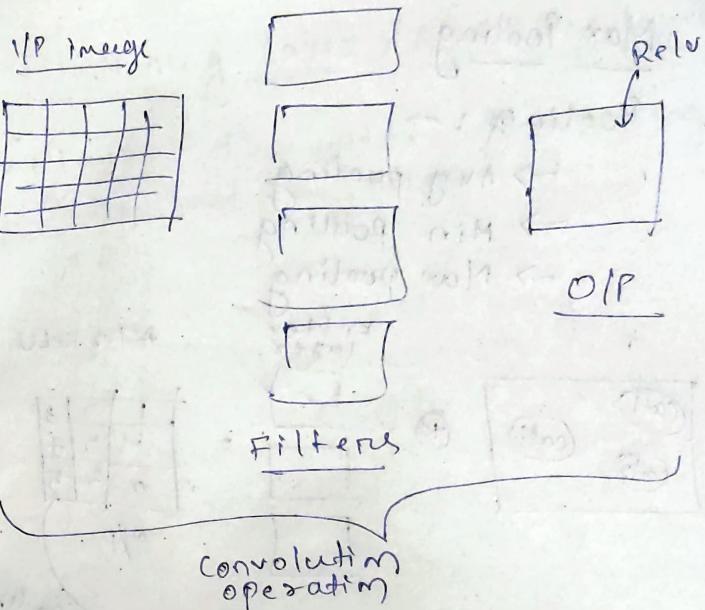
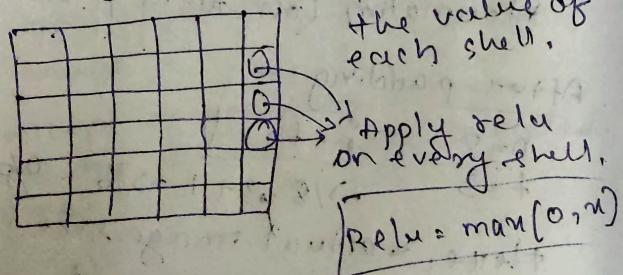
$$\begin{aligned}
 n &= 6 \text{ (Actual image size)} \\
 p &= 1 \text{ (padding layer)} \\
 f &= 3 \text{ (filter size)} \\
 s &= 1 \text{ (stride value)}
 \end{aligned}$$

Updated formula:

$$\begin{aligned}
 &\boxed{n + 2p - f + 1} \quad , \text{ if } s = 1 \\
 &= 6 + 2(1) - 3 + 1 \\
 &= 6 + 2 - 3 + 1 \\
 &= 6
 \end{aligned}$$

- As in ANN during the back propagation we need to update the weights.
- Like that in CNN we need to update the filters based on the input image.
- Initially we may initialize some filters randomly, but later we need to update the filters.

After convolution takes place, we get a 6×6 image, then apply relu & activation function on the value of each shell.

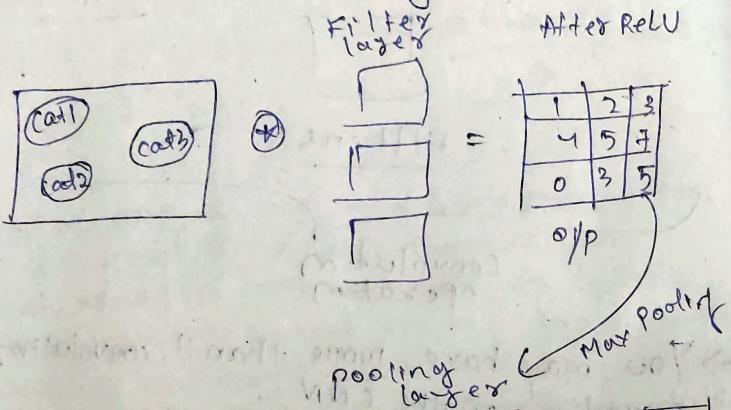


→ You may have more than 1 convolution operation on the CNN.

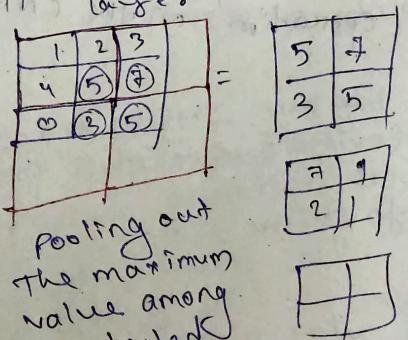
3 Max Pooling

Pooling :-

- Avg pooling
- Min pooling
- Max pooling



here, we're taking stride = 2 (by default)



Avg Pooling → In avg pooling, find out the avg from the selected values.

Min Pooling → Find out the minimum value among the selected.

4 Flattening Layer / ANN Dense layer

- Flatten the entire Max Pooling O/p. (Just elongating it)

5
7
3
5
7
9
2
1

Just attaching downwards in a single column.

Overall View

